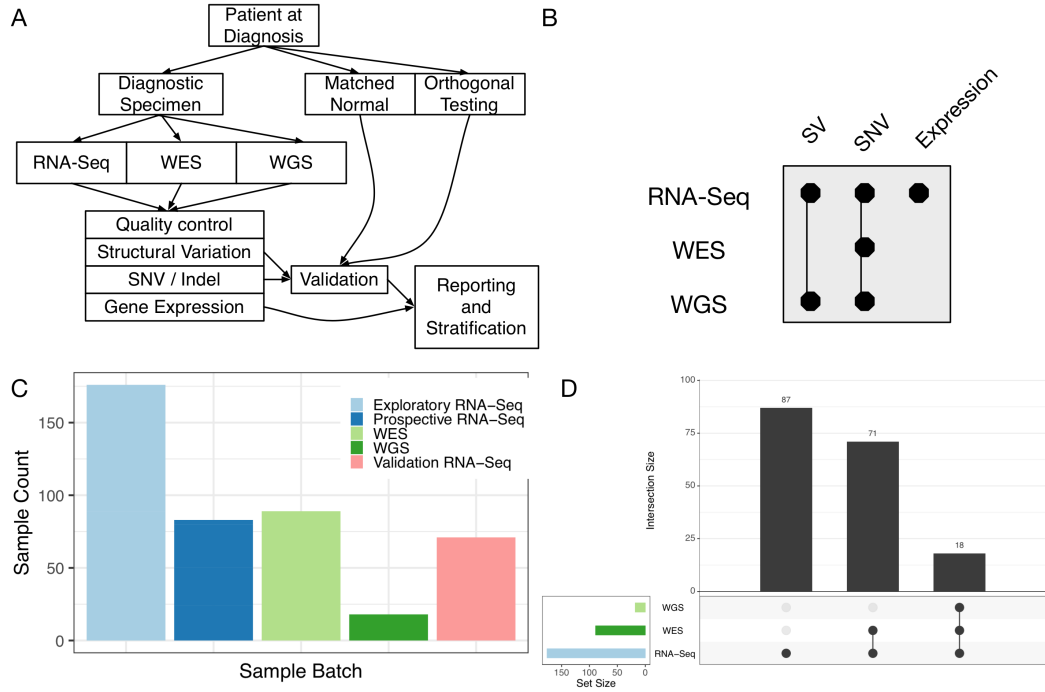


1 **Supplementary Information**

2

3 **Supplementary Figures**



4

5 **Supplementary Figure 1.** Experimental design overview. **(A)** Experimental workflow

6 for the exploratory RNA-Seq, WES, and WGS libraries. Diagnostic tumour specimens

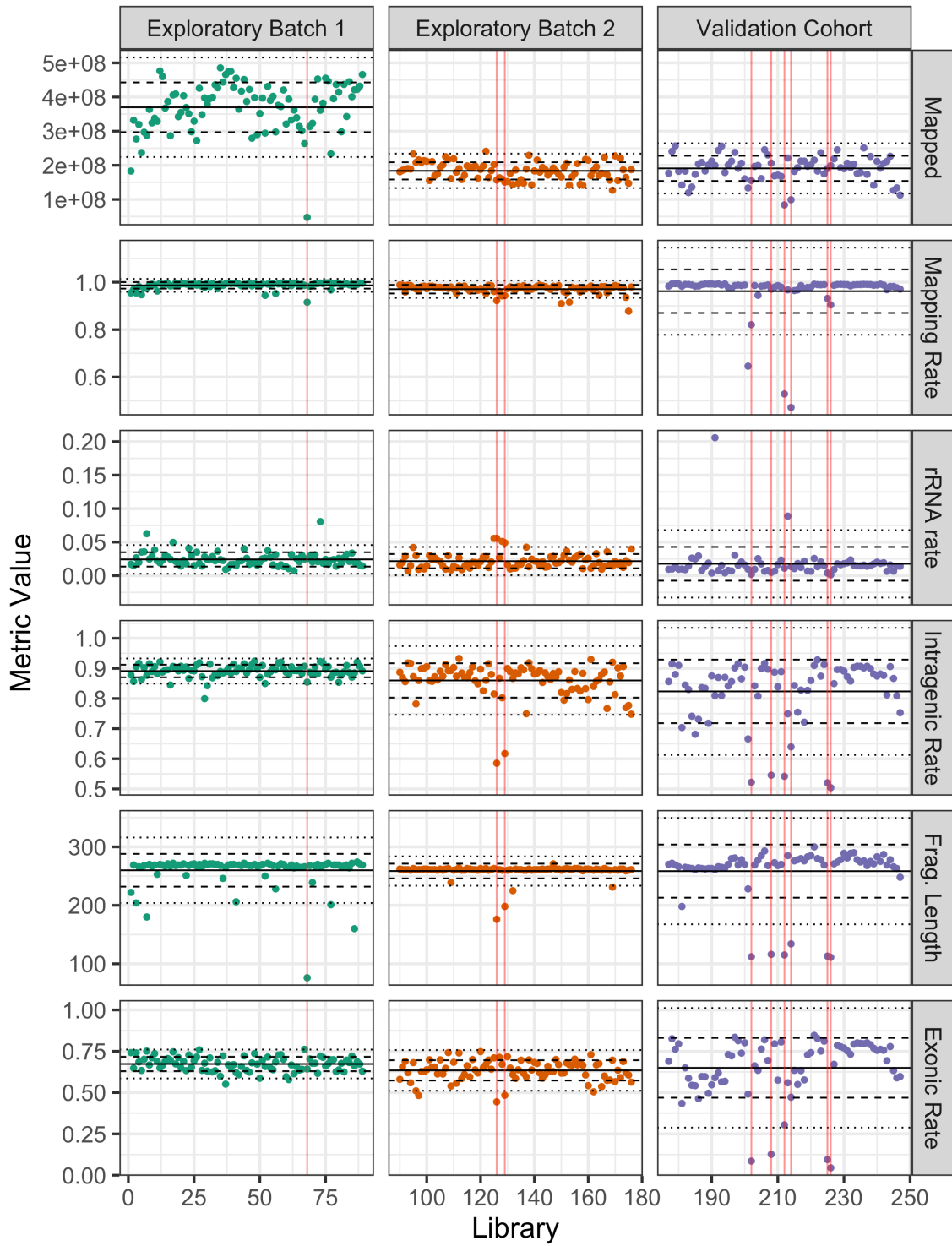
7 were obtained for each patient, with matched normal material used to validate selected

8 somatic SNV and short indel alterations. **(B)** Informatic analyses performed for each

9 sequencing platform. SV: structural variation, SNV: single nucleotide variant. **(C)**

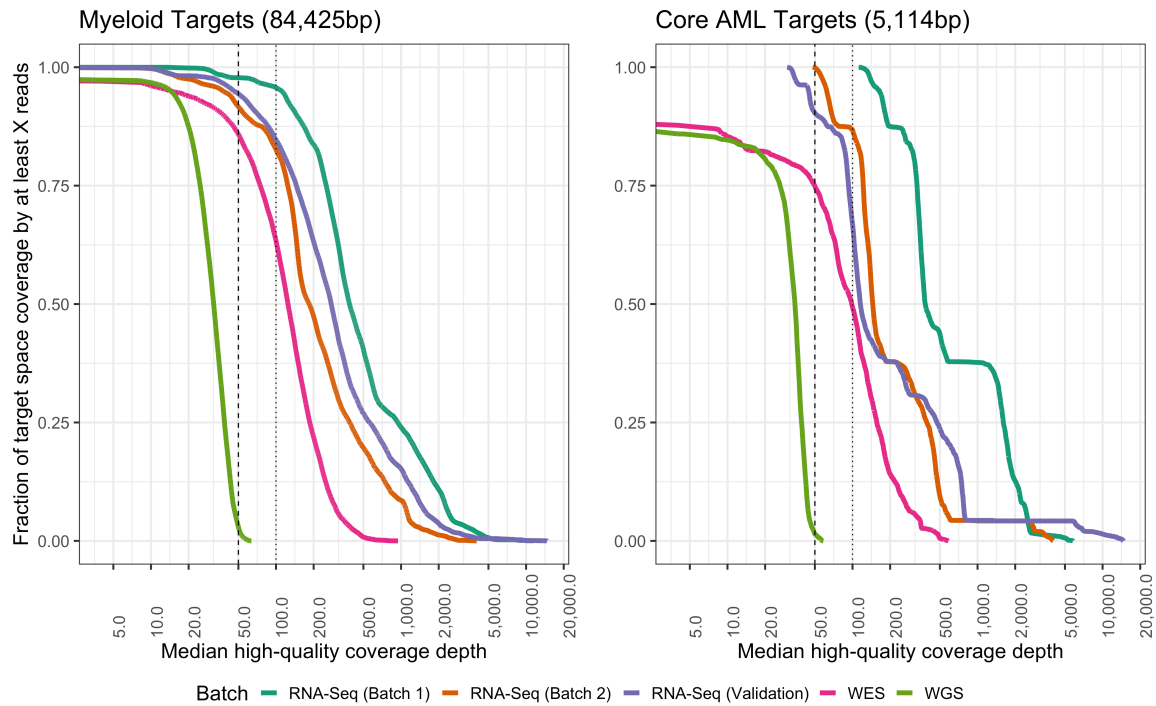
10 Sample count by project batch and sequencing platform. **(D)** Patient samples with

11 matched RNA-Seq, WES, and WGS libraries in the exploratory cohort.



14 **Supplementary Figure 2.** Levey-Jennings quality control charts for selected quality
15 metrics for RNA-Seq libraries from the AML PMP. In each panel, the x-axis represents
16 the sample number within the project, while y-axis represents a distinct quality metric.
17 The solid horizontal line represents the mean value for each metric within each
18 sequencing batch, while the dashed and dotted lines represent one and two standard
19 deviations away from the mean, respectively. ‘Mapped’: total number of mapped reads.
20 ‘Mapping Rate’: proportion of total reads which were mapped to the reference genome.
21 ‘rRNA rate’: proportion of reads originating from rRNAs. ‘Intragenic rate’: proportion of
22 reads that map within genes. ‘Fragment length’: mean observed fragment length. ‘Exonic
23 rate’: proportion of reads mapping within exons. Failed samples are indicated with red
24 vertical lines.

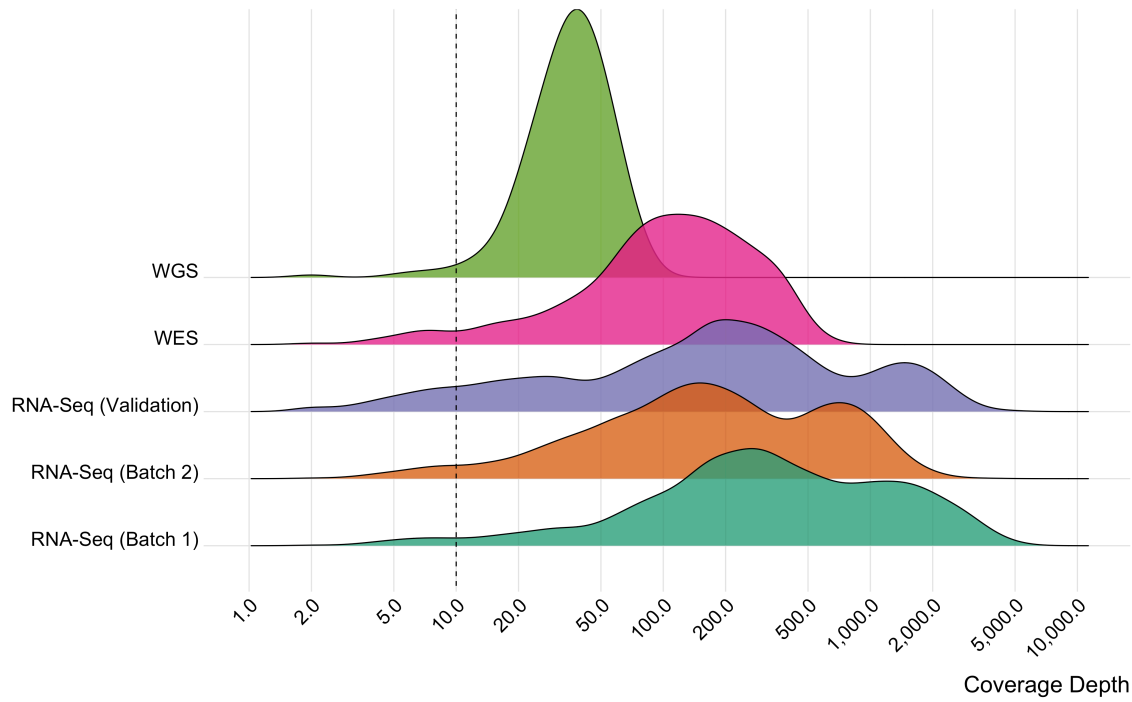
25



26

27 **Supplementary Figure 3.** Empirical distribution for sequence coverage depth. Myeloid
28 panel targets (left panel) and core clinical targets (*CEBPA*, *DNMT3A*, *FLT3*, *IDH1*,
29 *IDH2*, *KIT*, and *NPM1*) (right panel). The y-axis represents the proportion of bases
30 covered by at least the number of reads indicated by the x-axis. Vertical dashed lines
31 indicate coverage thresholds of 50x and 100x. Sequencing platform and study batch are
32 indicated by line colour.

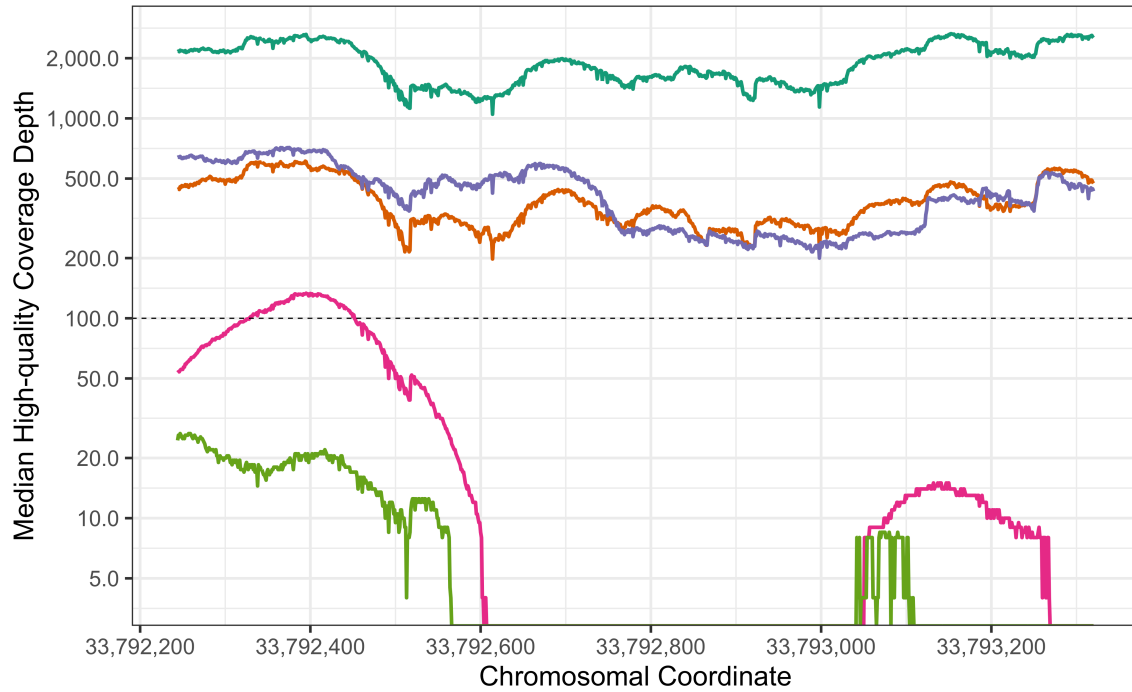
33



34

35 **Supplementary Figure 4.** Observed sequencing coverage depth for called variants across
36 sequencing platforms and batches in the AML PMP cohort. Sample groups are coloured
37 as in Supplementary Figure 3.

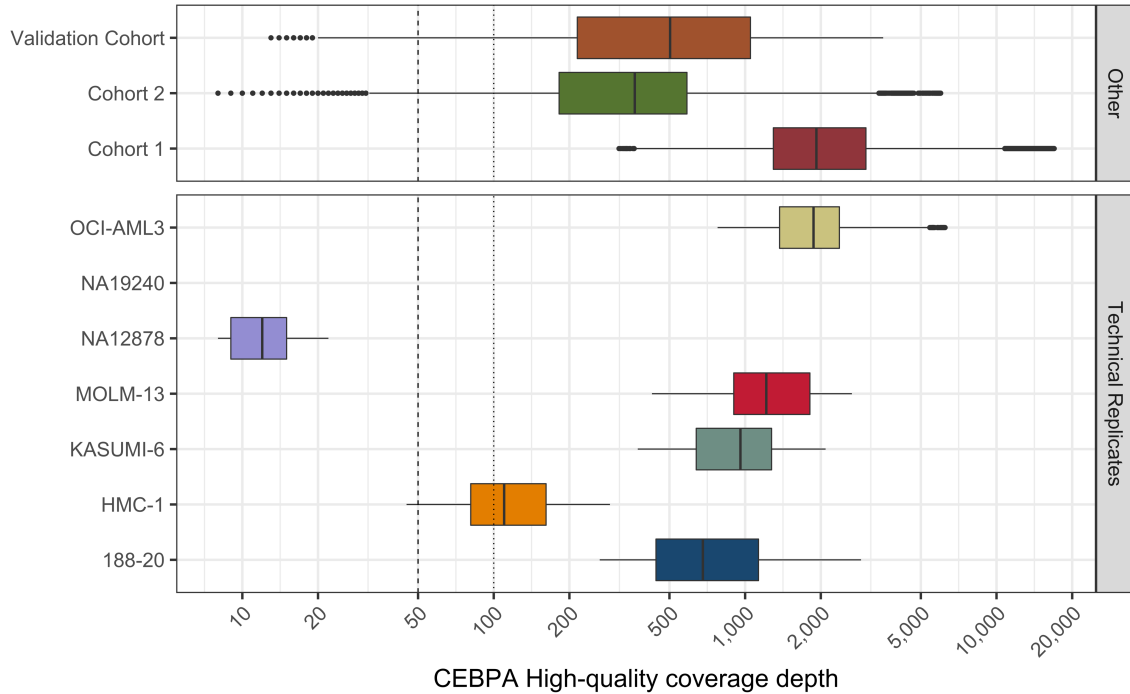
38



39

40 **Supplementary Figure 5.** Median high-quality sequence coverage depth over *CEBPA* in
41 the AML PMP cohort. Each point represents the median coverage across all samples for a
42 particular sequencing platform and batch. The 'Chromosomal Coordinate' indicates the
43 hg19 position along chromosome 19. The line colour indicates the sequencing platform
44 and study batch, as in Supplementary Figure 3.

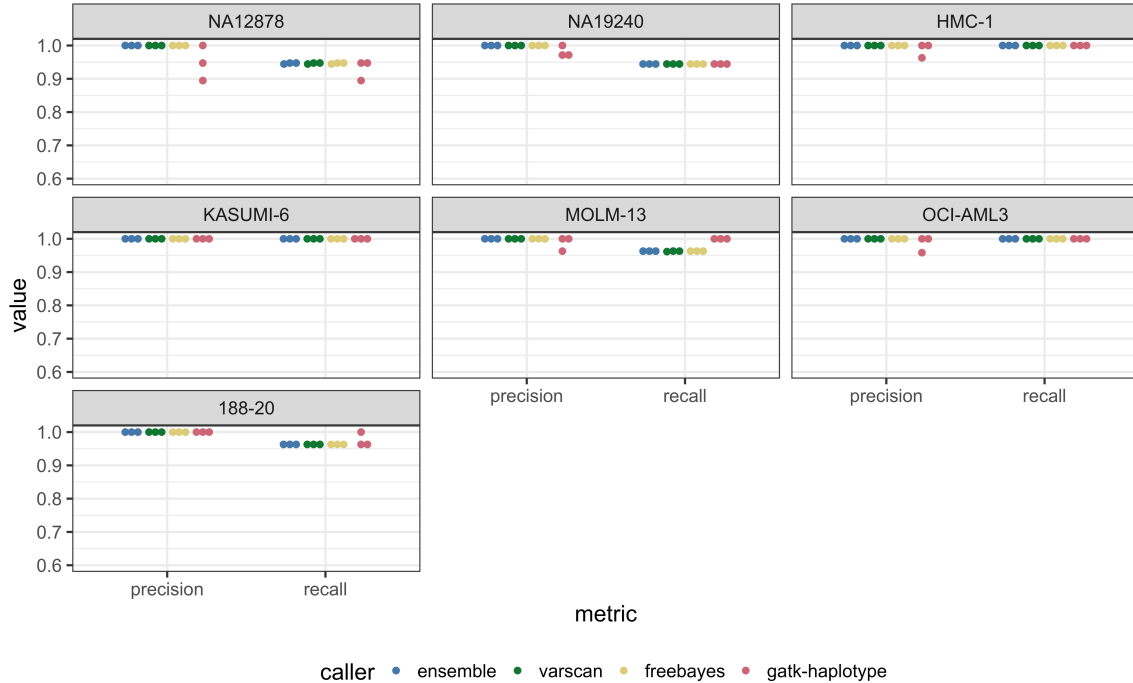
45



46

47 **Supplementary Figure 6.** High-quality coverage depth in *CEBPA* for technical replicate
48 samples compared to other RNA-Seq and ssRNA-Seq samples. For each sample, high-
49 quality sequence coverage depth at 1,077 sites in *CEBPA* were assessed. Each box plot
50 represents pooled data across retained samples from Cohort 1 (n = 88), Cohort 2 (n = 85),
51 and the RNA-Seq Validation Cohort (n = 28). For cell lines OCI-AML3, NA19240,
52 NA12878, MOLM-13, KASUMI-6, HMC-1, and patient sample 188-20, three technical
53 replicates per condition are shown. For all box plots, the median is shown with a vertical
54 line, with the edges of the box corresponding to the first and third quartiles of the
55 distribution, with whiskers extending to 1.5 * the inter-quartile range of the distribution.

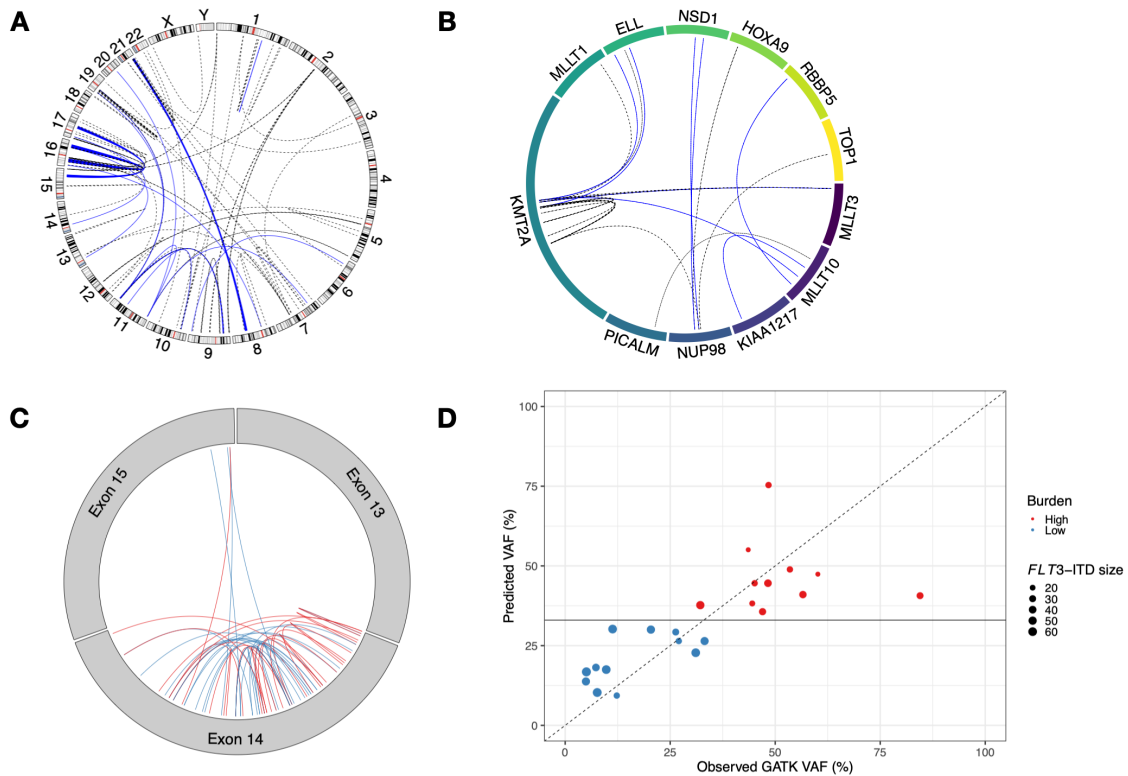
56



57

58 **Supplementary Figure 7.** Precision and recall by sample and variant caller for replicated
59 RNA-Seq libraries from the AML PMP validation cohort. Precision: $TP / (TP + FP)$,
60 Recall: $TP / (TP + FN)$. Libraries which failed sample QC were excluded.

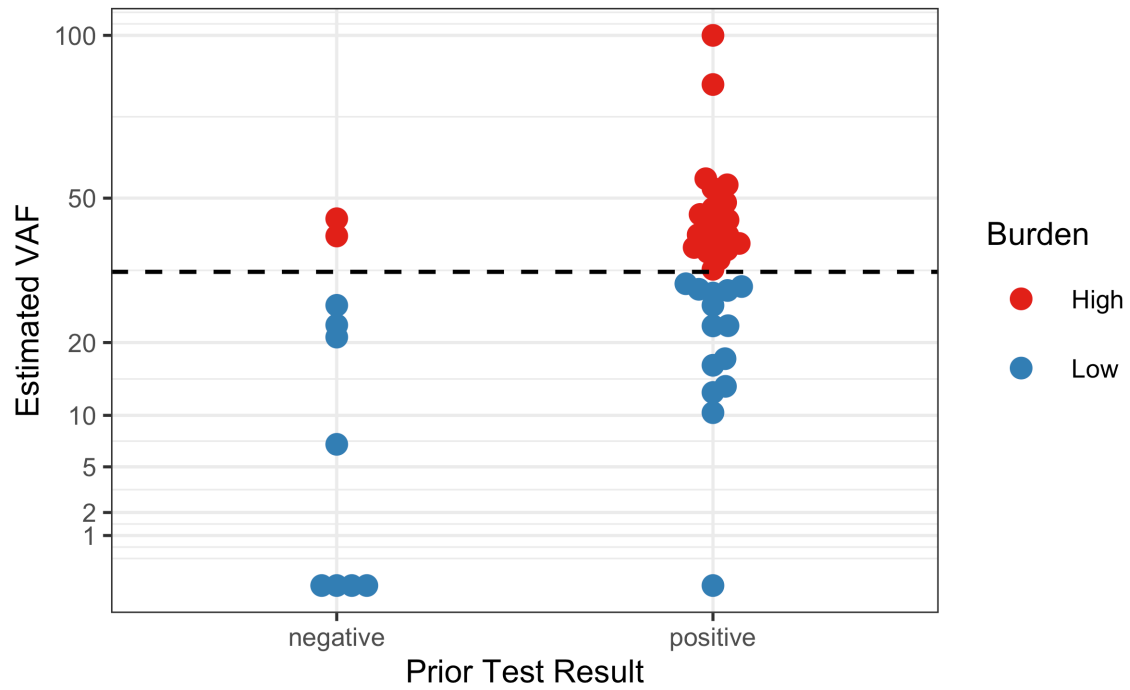
61



62

63 **Supplementary Figure 8.** RNA-Seq structural variation in the AML PMP exploratory
64 cohorts. **(A)** Filtered gene fusion events detected by RNA-Seq. Each arc represents a
65 distinct set of fusion partners, for known (blue) and novel (black) rearrangements. **(B)**
66 Structural rearrangements involving *KMT2A*-related genes. Each gene (except *KMT2A*) is
67 scaled to the same size, for known (blue) and novel (black) events. **(C)** Intra-gene
68 structural variation in *FLT3*. Each circle segment represents a single exon of *FLT3*, and
69 the zoomed panels show the coordinates of internal tandem duplication breakpoints. **(D)**
70 Predicted VAF for *FLT3*-ITD events based on linear modelling, compared to the
71 observed VAF from GATK HaplotypeCaller for *FLT3*-ITD events detected by both
72 GATK and trans-ABYSS.

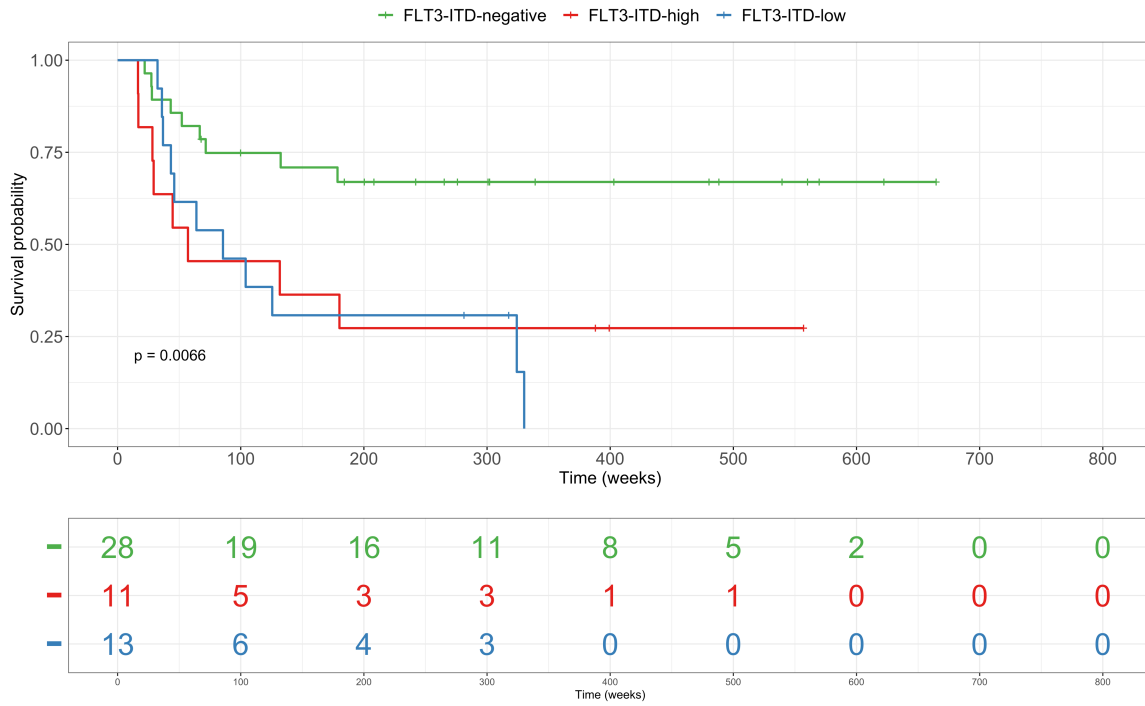
73



74

75 **Supplementary Figure 9.** Estimated VAF for *FLT3*-ITD events detected or not detected
76 by prior clinical assay in the AML PMP exploratory cohorts. The horizontal dashed line
77 indicates the estimated VAF of 33% used to discriminate high-burden from low-burden
78 *FLT3*-ITD events.

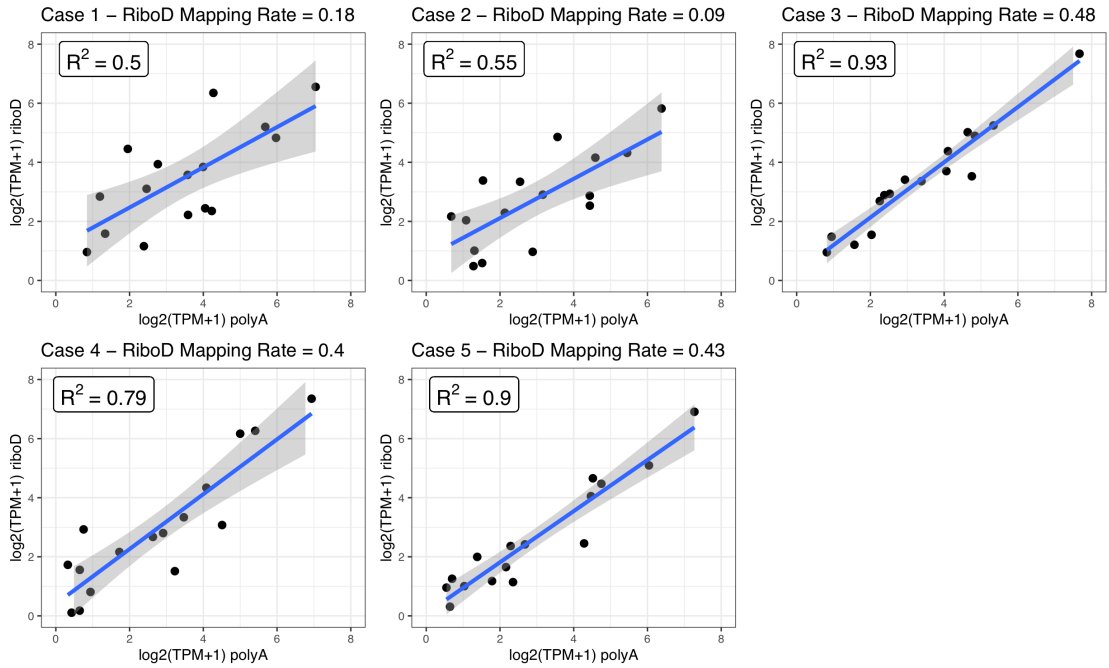
79



80

81 **Supplementary Figure 10.** Overall survival for normal-karyotype AML patients in the
 82 AML PMP exploratory cohorts by *FLT3*-ITD burden, compared to patients with no
 83 *FLT3*-ITD event.

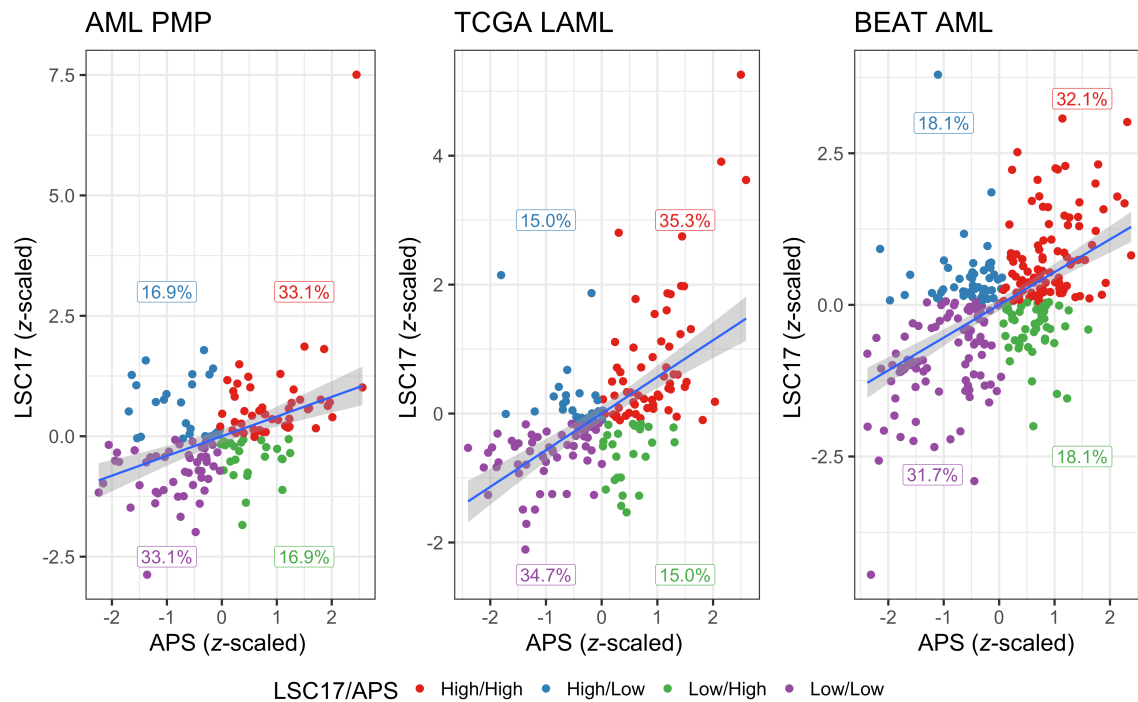
84



85

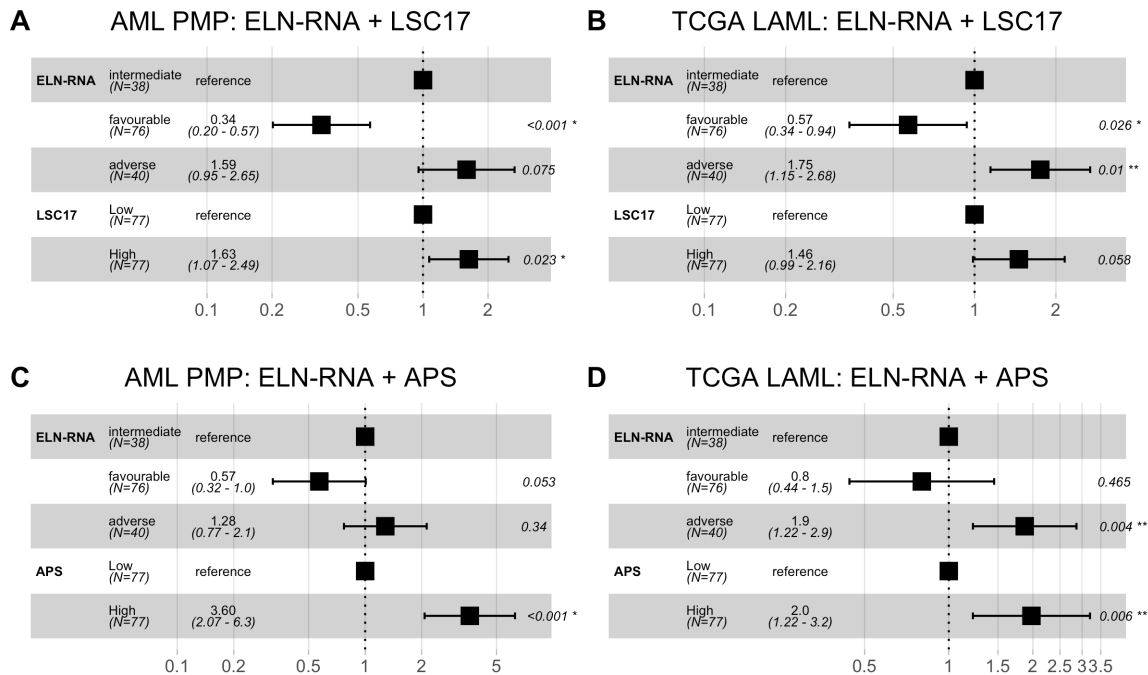
86 **Supplementary Figure 11.** Comparison between matched RNA-Seq libraries prepared
87 with polyA or ribodepletion protocols. In each panel, the mapping rate for the
88 ribodepleted library is indicated in the panel title. Each panel also indicates a linear
89 regression of the two variables, with 95% confidence interval.

90



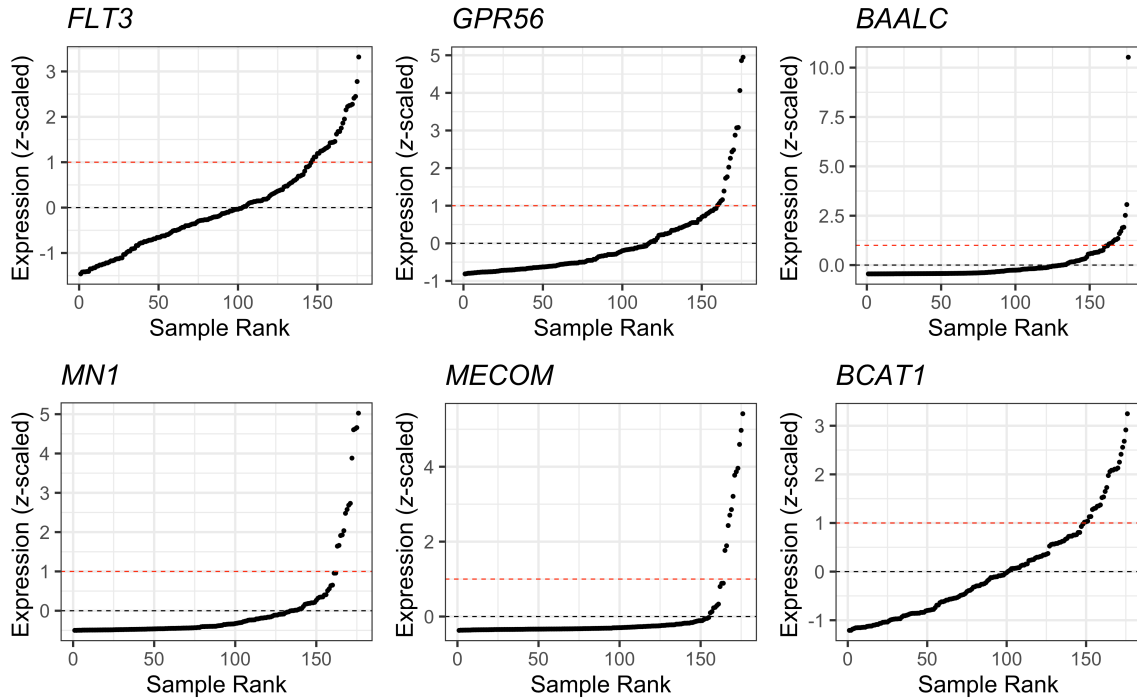
91

92 **Supplementary Figure 12.** Scaled APS and LSC17 values in the AML PMP, TCGA
93 LAML, and BEAT AML cohorts. For each cohort, APS and LSC17 values were
94 standardized by calculating z scores. High/low categories used for patient stratification
95 are indicated by colour. Each panel also indicates a linear regression of the two variables
96 (indicated with a blue line), with 95% confidence interval.



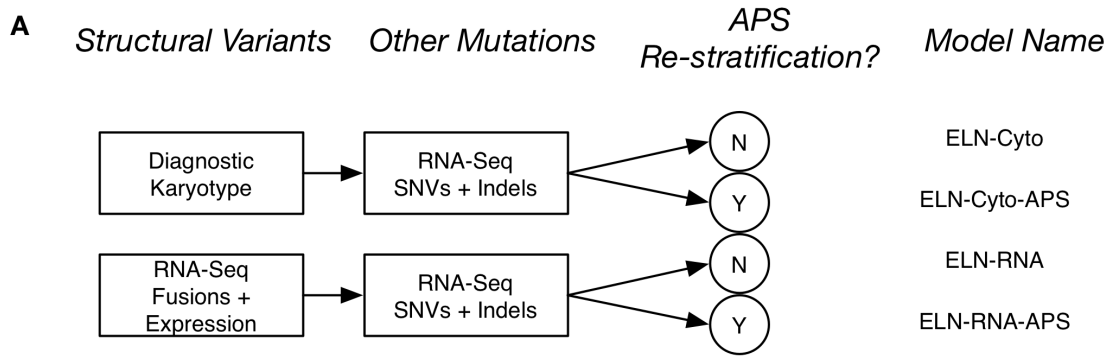
98

99 **Supplementary Figure 13.** Forest plots for multivariate survival models based on ELN-
 100 RNA stratifications with either the LSC17 or APS gene signatures. Multivariate Cox
 101 proportional hazards regression models were fit, modeling survival as a function of the
 102 ELN-RNA stratification and an expression score (either LSC17 or APS). In each panel,
 103 the hazard ratio is indicated on the x-axis, and model variables are indicated on the y-
 104 axis. In each model, ‘Intermediate-risk’ and ‘Low’ values for the gene signatures were set
 105 as the reference levels. In each panel, the estimated hazard ratio for a given variable is
 106 indicated by a square, with 95% confidence intervals. *p*-values associated with each
 107 variable are indicated at the right of each plot.



109

110 **Supplementary Figure 14.** Single-gene expression outliers. Samples are ranked by gene
 111 expression, measured as a z score (where $z = (x - \mu)/\sigma$). Black dashed line = mean
 112 expression, red dashed line = high outlier cutoff.



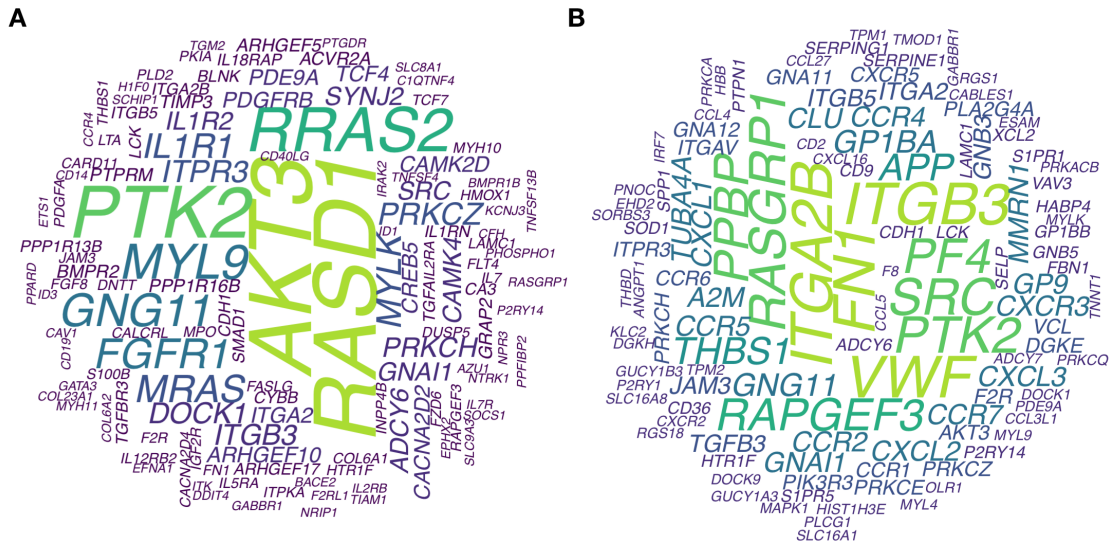
B

ELN-Cyto	ELN-RNA	n
Favourable	Favourable	74
Intermediate	Favourable	1
Intermediate	Intermediate	25
Intermediate	Adverse	7
Adverse	Favourable	1
Adverse	Intermediate	13
Adverse	Adverse	33
Discrepant		22

C

ELN-Cyto-APS	ELN-RNA-APS	n
Favourable	Favourable	77
Intermediate	Intermediate	6
Intermediate	Adverse	2
Adverse	Intermediate	7
Adverse	Adverse	62
Discrepant		9

115 **Supplementary Figure 15.** Comparison of patient stratification models. (A) Schematic
 116 overview of stratification models applied. (B) Comparison of patient stratifications
 117 between the ELN-Cyto and ELN-RNA models. (C) Comparison of stratifications
 118 between the ELN-Cyto-APS and ELN-RNA-APS models.



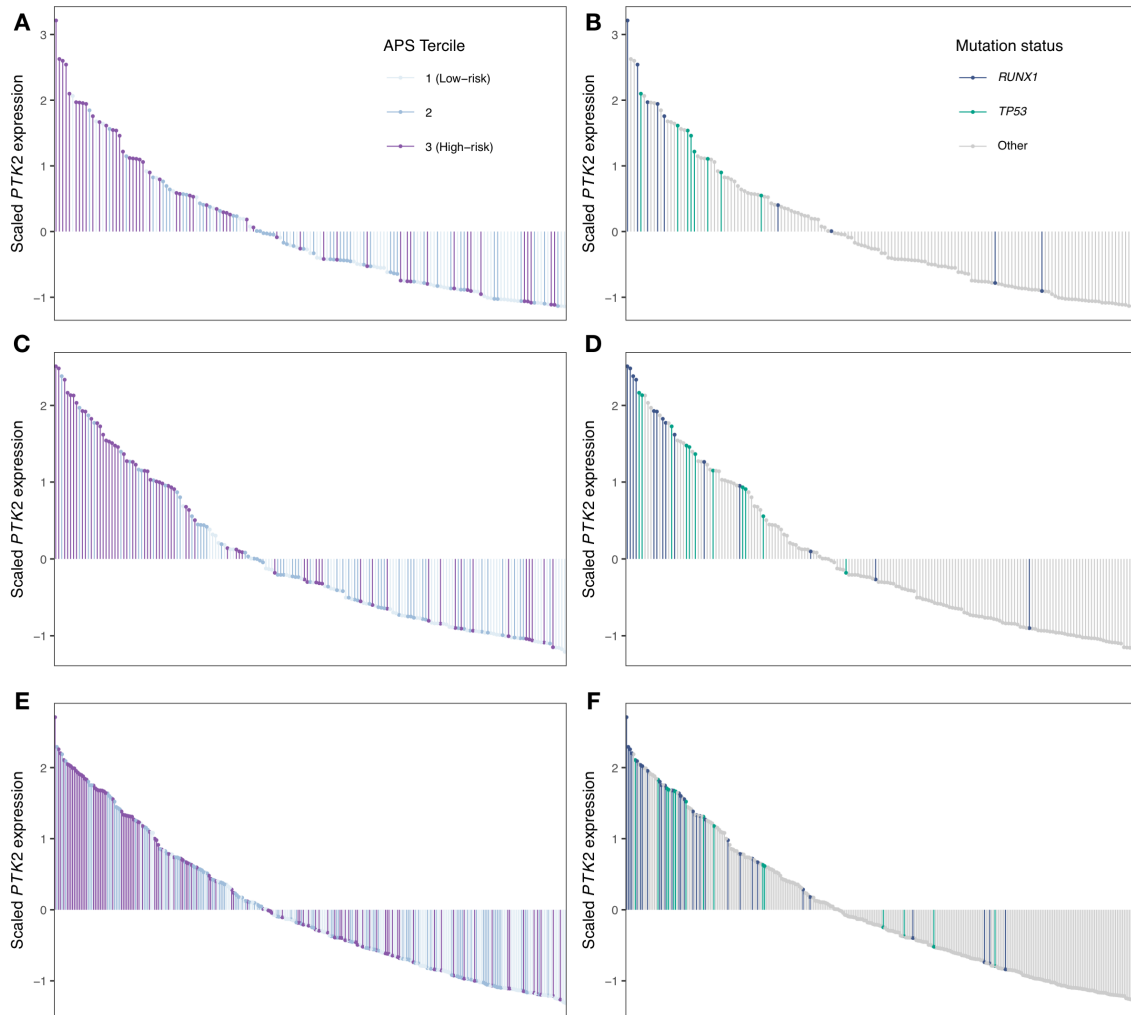
120

121 **Supplementary Figure 16.** Recurrently used molecules in enriched pathways from IPA

122 (A) and GSEA (B) pathway enrichment analyses. For each analysis, IPA molecules or

123 leading edge genes from enriched pathways were extracted and summarized, and

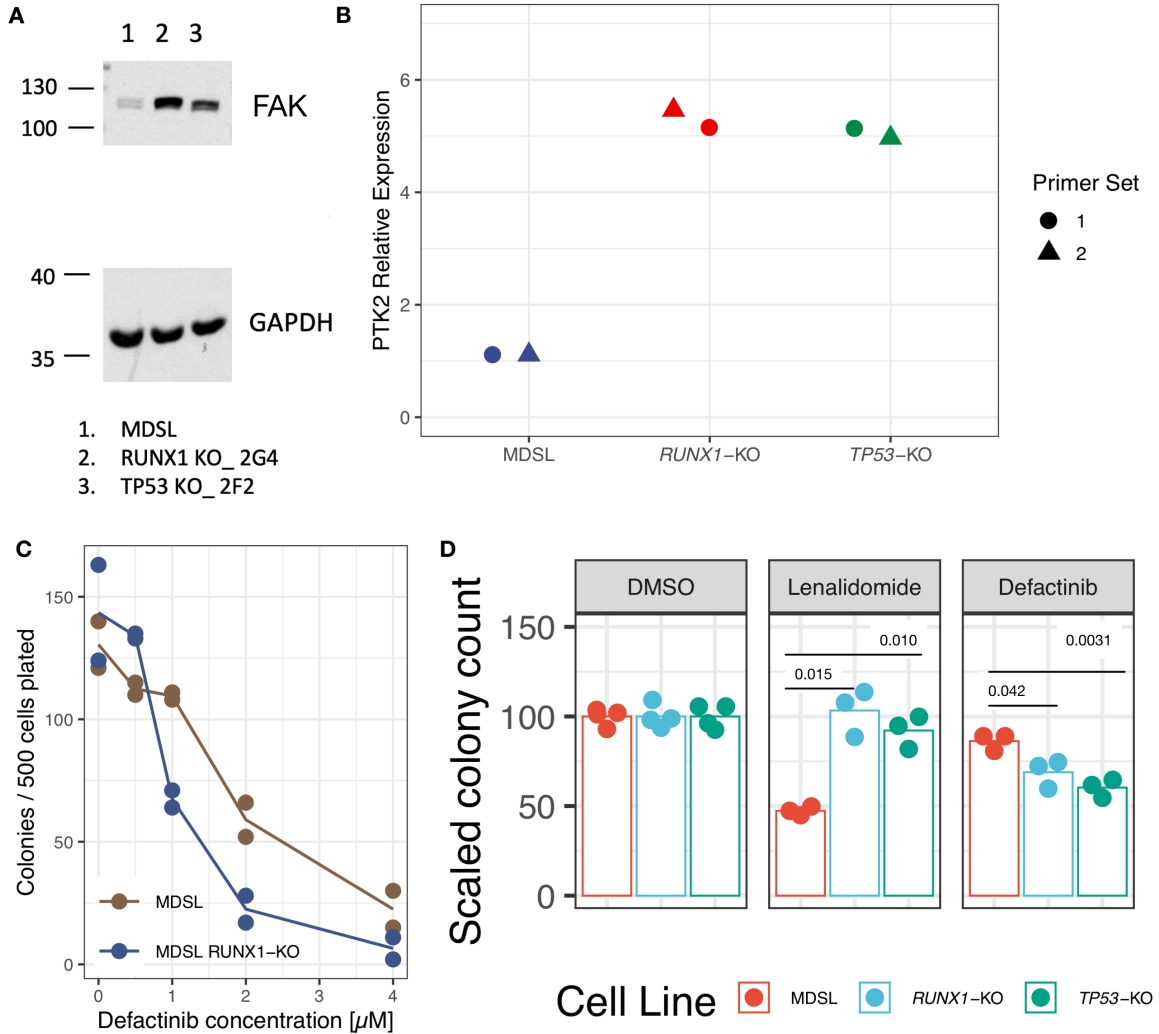
124 visualized as word clouds.



126

127 **Supplementary Figure 17.** Patients ranked by *PTK2* expression. (A, C, E) Patients are
 128 ranked by *PTK2* expression, and coloured by APS tercile for the AML PMP (A), TCGA
 129 LAML (C), and BEAT AML (E) cohorts. (B, D, F) Patients are ranked by *PTK2*
 130 expression, and coloured mutation status for the AML PMP (B), TCGA LAML (D), and
 131 BEAT AML (F) cohorts.

132



133

134 **Supplementary Figure 18.** FAK inhibition in MDSL cell line derivatives. (A) Western
 135 blot for FAK protein expression in MDSL cell lines with *RUNX1* or *TP53* CRISPR
 136 knockout, using cell lines generated by Martinez-Hoyer *et al.*¹ Each blot was repeated
 137 twice, with similar results. (B) qPCR relative expression for *PTK2*. (C) Colony forming
 138 cell-count dose-response curve for defactinib. (D) Colony forming cell count assays for
 139 cells treated with DMSO, 1 μ m Lenalidomide, or 1 μ m Defactinib, with two-sided *t*-test *p*
 140 values indicated.

141 **Supplementary References**

- 142 1. Martinez-Høyer, S. *et al.* Loss of lenalidomide-induced megakaryocytic
143 differentiation leads to therapy resistance in del(5q) myelodysplastic syndrome.
144 *Nature cell biology* (2020). doi:[10.1038/s41556-020-0497-9](https://doi.org/10.1038/s41556-020-0497-9)